

I Love Machine Learning - Separate Token

Start Token

[CLS]

I

Love

machine

Learning [Sep]

Token Embedding

Position Embedding

Input Embedding

T_0

T_1

T_2

T_3

T_4

T_5

P_0

P_1

P_2

P_3

P_4

P_5

768x1 - Best BERT
1536x1 - Best Long

Embed. vectors

768x1
or 1536x1

Capture sequence info

$T_0 + P_0$

$T_1 + P_1$

$T_2 + P_2$

$T_3 + P_3$

$T_4 + P_4$

$T_5 + P_5$

Self-Attention

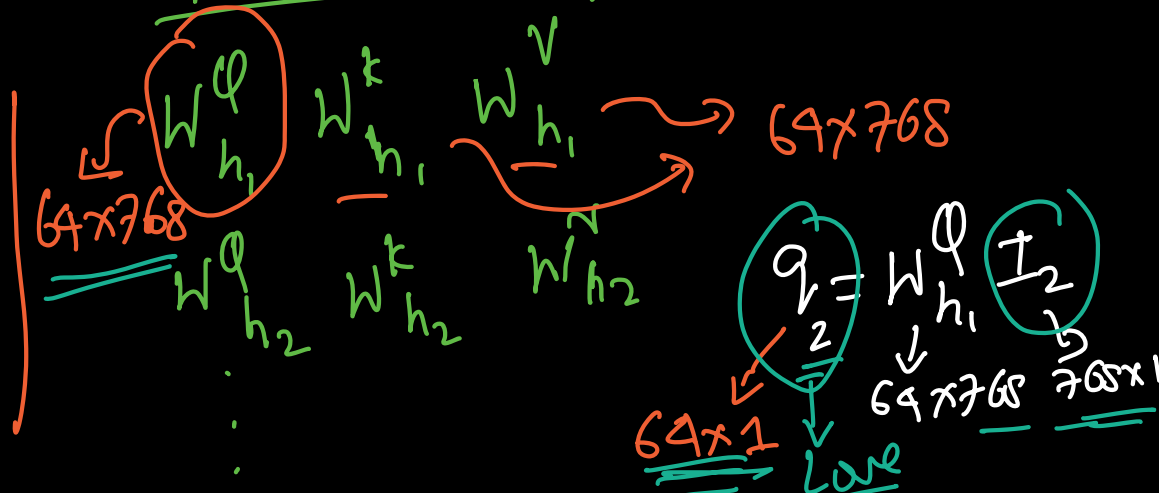
W^Q, W^K, W^V

Weight Matrices for Query, Key, Value

Q - query
 K - key
 V - value

Multi-Head self-Attention

Assume
12 heads
Head Embed
dim = $\frac{768}{12} = 64$

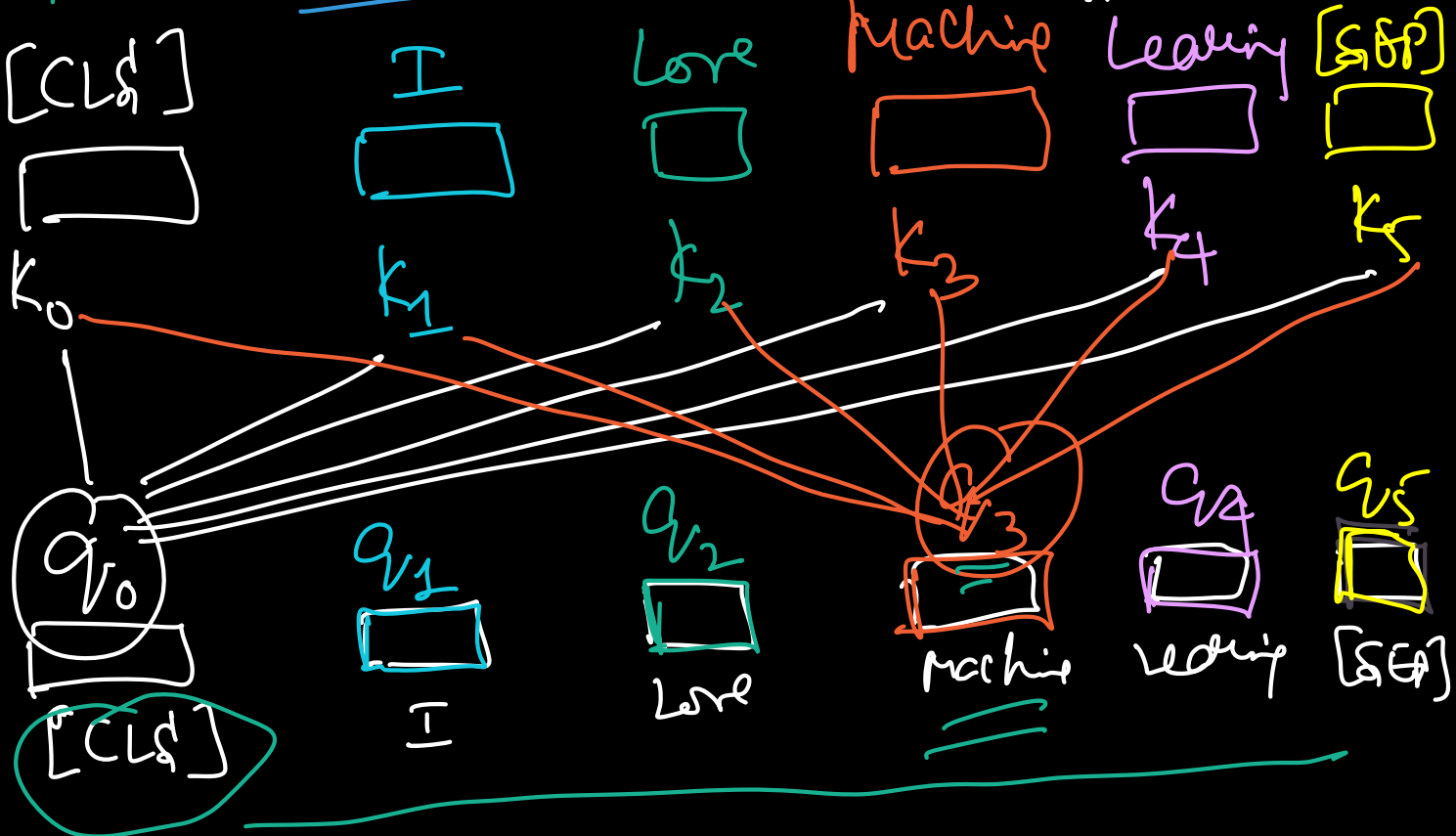


$$I_0 = P_0 + T_0, \quad I_1 = T_1 + P_1, \quad \dots$$



Pick one head

$$\left. \begin{aligned} q_0 &= W^Q I_0 \\ k_0 &= W^K I_0 \\ v_0 &= W^V I_0 \end{aligned} \right\}$$



$$\begin{matrix}
 & & 0 & 1 & 2 & 3 & 4 & 5 \\
 \begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix}
 \underline{q_{0 \cdot k_0}} & \underline{q_{0 \cdot k_1}} & \underline{q_{0 \cdot k_2}} & q_{0 \cdot k_3} & q_{0 \cdot k_4} & q_{0 \cdot k_5} \\
 & & & & & \\
 & & & & & \\
 q_{3 \cdot k_0} & q_{3 \cdot k_1} & q_{3 \cdot k_2} & q_{3 \cdot k_3} & q_{3 \cdot k_4} & q_{3 \cdot k_5} \\
 & & & & & \\
 & & & & &
 \end{bmatrix} & \left. \vphantom{\begin{matrix} 0 \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix}} \right\} \\
 \end{matrix}$$

6x6

Attention is captured by "dot product"

In general $A_{n \times n}$ matrix

what is n?
 $n = \# \text{tokens}$

Let

$$Q = \begin{bmatrix} q_{v_0} \\ q_{v_1} \\ q_{v_2} \\ \vdots \\ q_{v_6} \end{bmatrix} \quad 6 \times 64$$

$$K = \begin{bmatrix} k_0 \\ k_1 \\ k_2 \\ \vdots \\ k_6 \end{bmatrix}$$

Is q_{v_0} a row vector or column vector?

Notice that:

$$\text{Attention (Un-Scaled)} = \begin{bmatrix} q_0 & \dots & q_{n-1} \\ \vdots & & \vdots \\ q_0 & \dots & q_{n-1} \end{bmatrix} \begin{bmatrix} k_0^T & k_1^T & \dots & k_{n-1}^T \end{bmatrix}$$

Q K^T

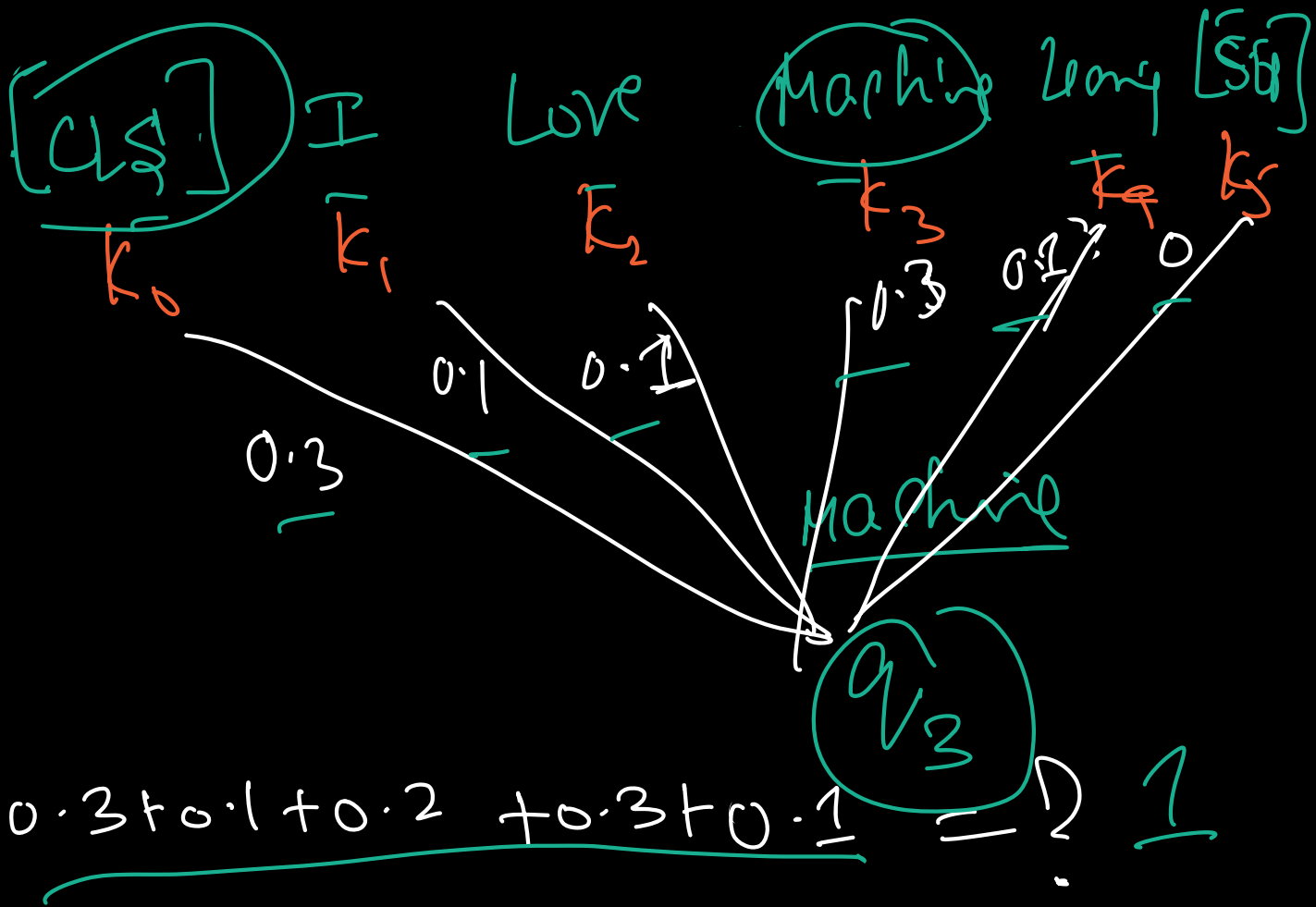
row vector column vector

Dot product $q_0 \cdot k_0 = \begin{bmatrix} \quad \quad \quad \end{bmatrix} \begin{bmatrix} \quad \quad \quad \end{bmatrix}$

Scaled dot-product self-attention

$$A = \frac{QK^T}{\sqrt{d}} = QK^T / \sqrt{64} = \underline{\underline{QK^T / 8}}$$

Convert $A \rightarrow$ Probabilistic Attention



Softmax

$$d_{ij} = \frac{e^{A_{ij}}}{\sum_j e^{A_{ij}}}$$

Scaled Dot-Product

Self-Attention
(Probabilistic)

Each row is a

prob. vector

$$Q = \left[\frac{e^{A_{ij}}}{\sum_j e^{A_{ij}}} \right]$$

Interpretation:-

How much should machine
place emphasis on the 6
tokens [SEP] I love machine learning [CLS]

It is decided by the vectors
in α .

Which vector in α fits
machine?

Query, key, values

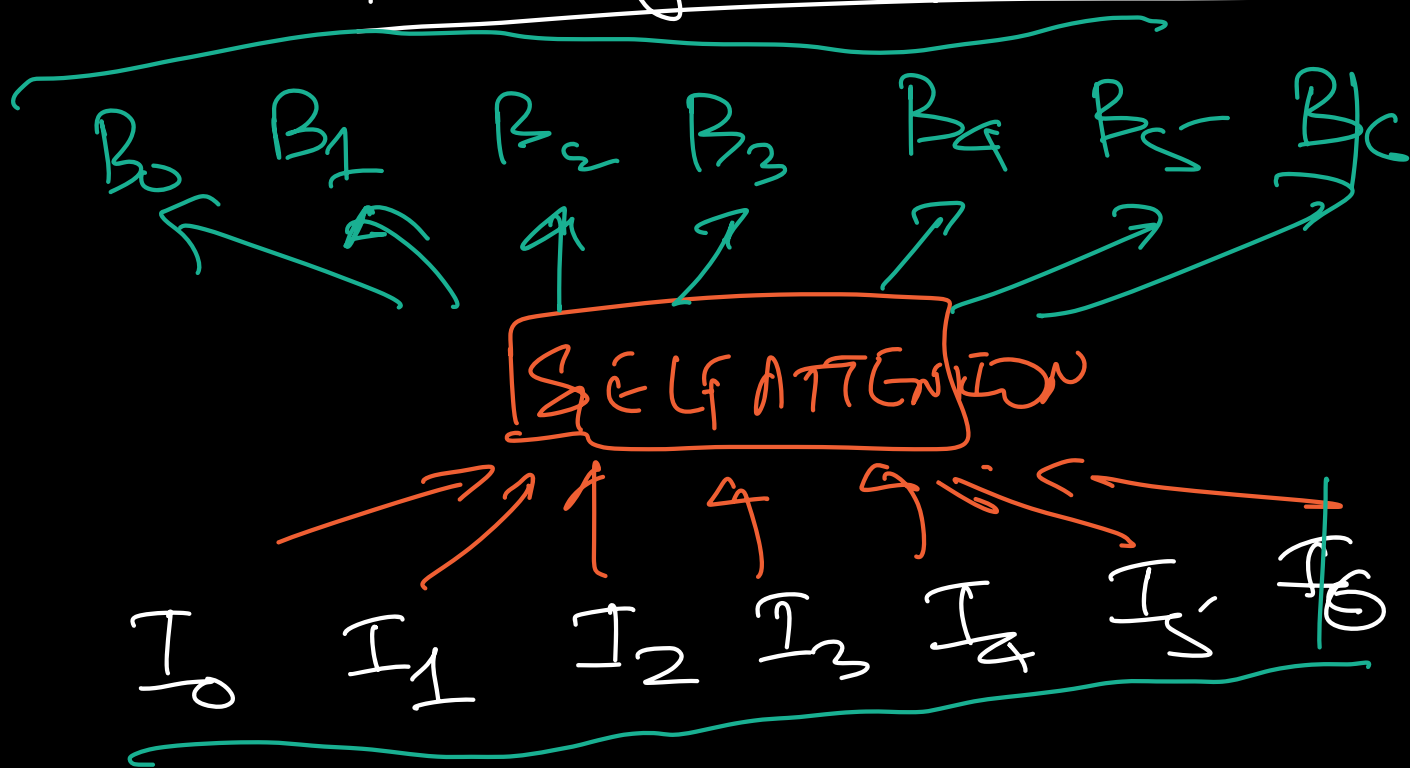
d_3 v-indexed
 d_4 i-index

Query - what you query
about

Key - what you match
with

Value - New understanding
of each token.

output of self-attention

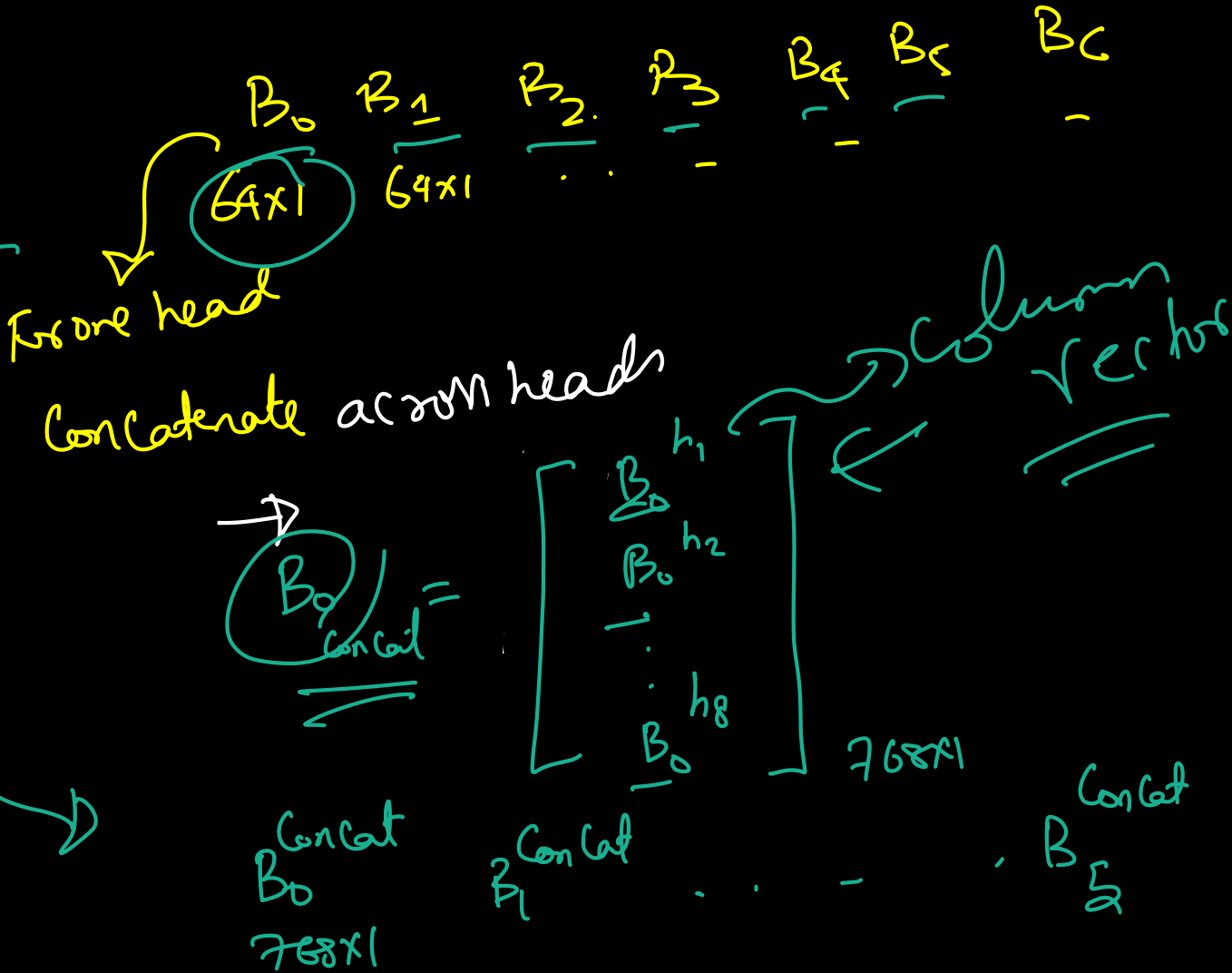


$B_i \rightarrow$ Blended understanding of all tokens for token i

So B_3 [Machine] = $d_{30} V_{[CLS]}$
 $+ d_{31} V_I + d_{32} V_{Love}$
 $+ d_{33} V_{machine} + d_{34} V_{leary}$
 $+ d_{35} V_{[SEP]}$

$\rightarrow V_{[CLS]} = W^V I_{[CLS]}$

Concatenate

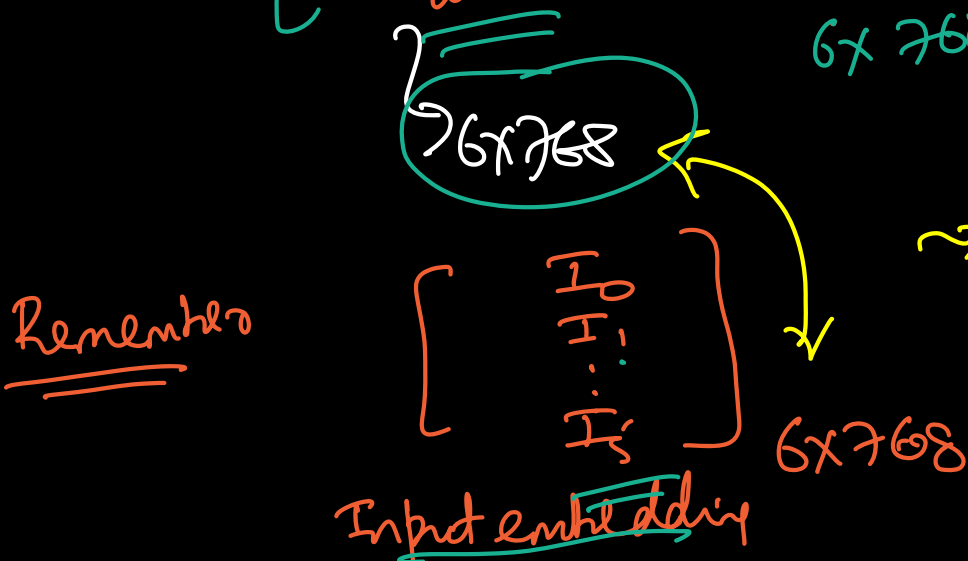


output projection

$$Z_{\text{attention}} = B_{\text{concat}}^T W_0$$

6x768 768x768

output projection



Input & output dim of embeddings out of the self-attention is the same!

We are not done yet!

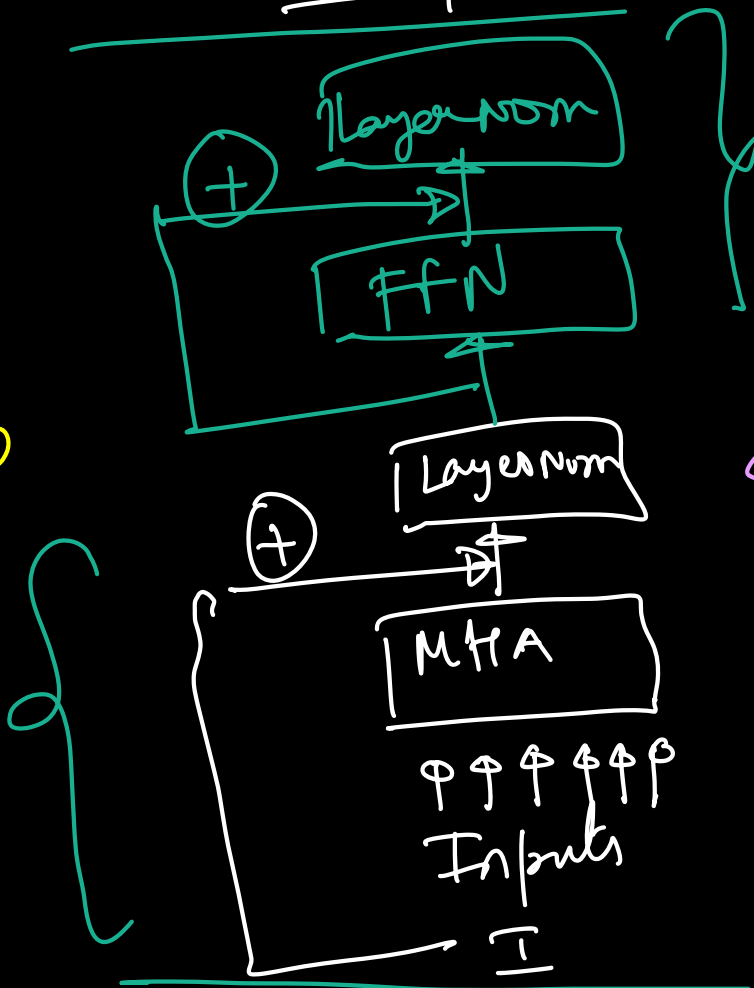
$$\tilde{X} = \text{Layer Norm} \left(\underline{I^{(0)}} + \text{Z attn} \right)$$

For Training stability

Adding residual connection ensures prev. embedding is still relevant!

Big picture

Layer 0



We just went through this component of one layer of BERT Transformer!

FFN

→ Each token undergoes a non-linear transformation independent of other tokens!!

1. Linear layer expansion

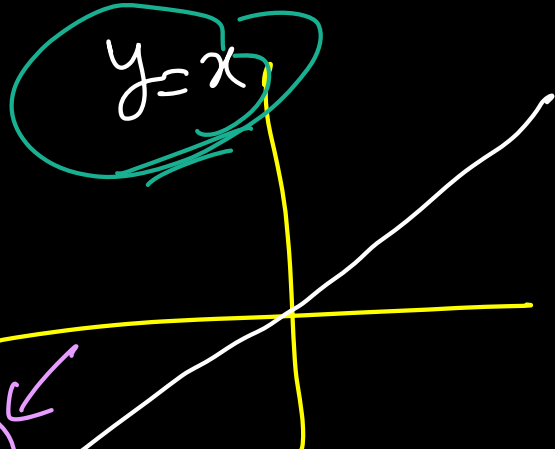
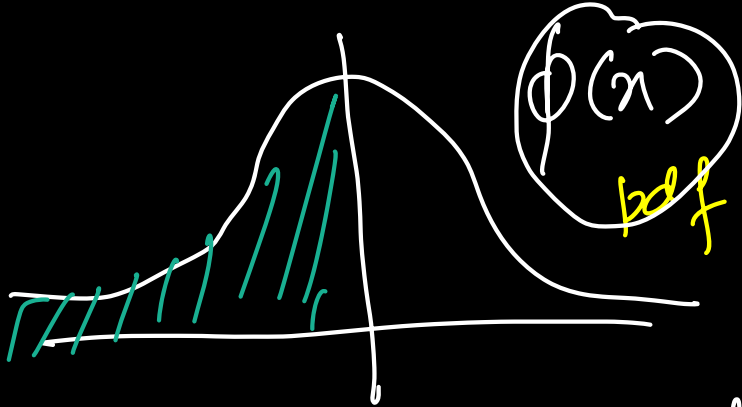
$$H = \sum W_1 + b_1$$

MHA → 6×768 768×3072 3072 → Dimension? 6×3072
↳ expansion

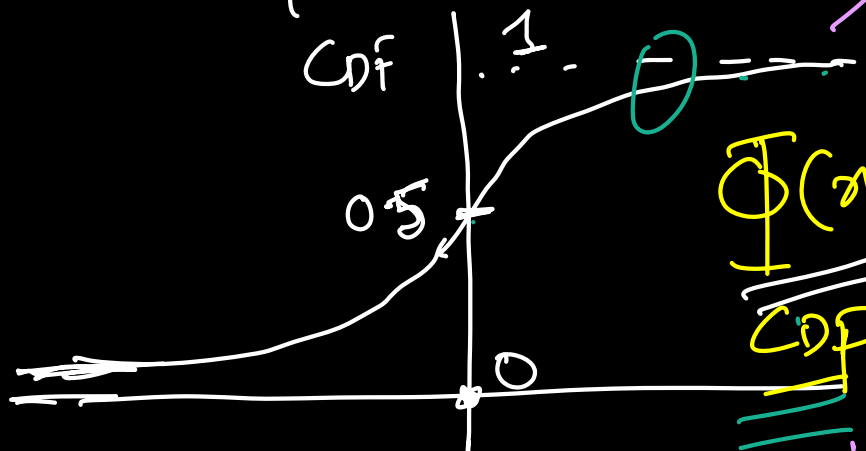
2. GELU activation → key NON-linearity (deep understanding)

$$\text{GELU}(x) = x \cdot \Phi(x)$$

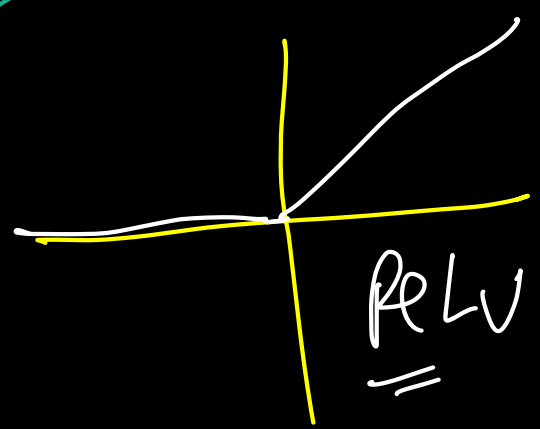
↳ CDF of Gaussian



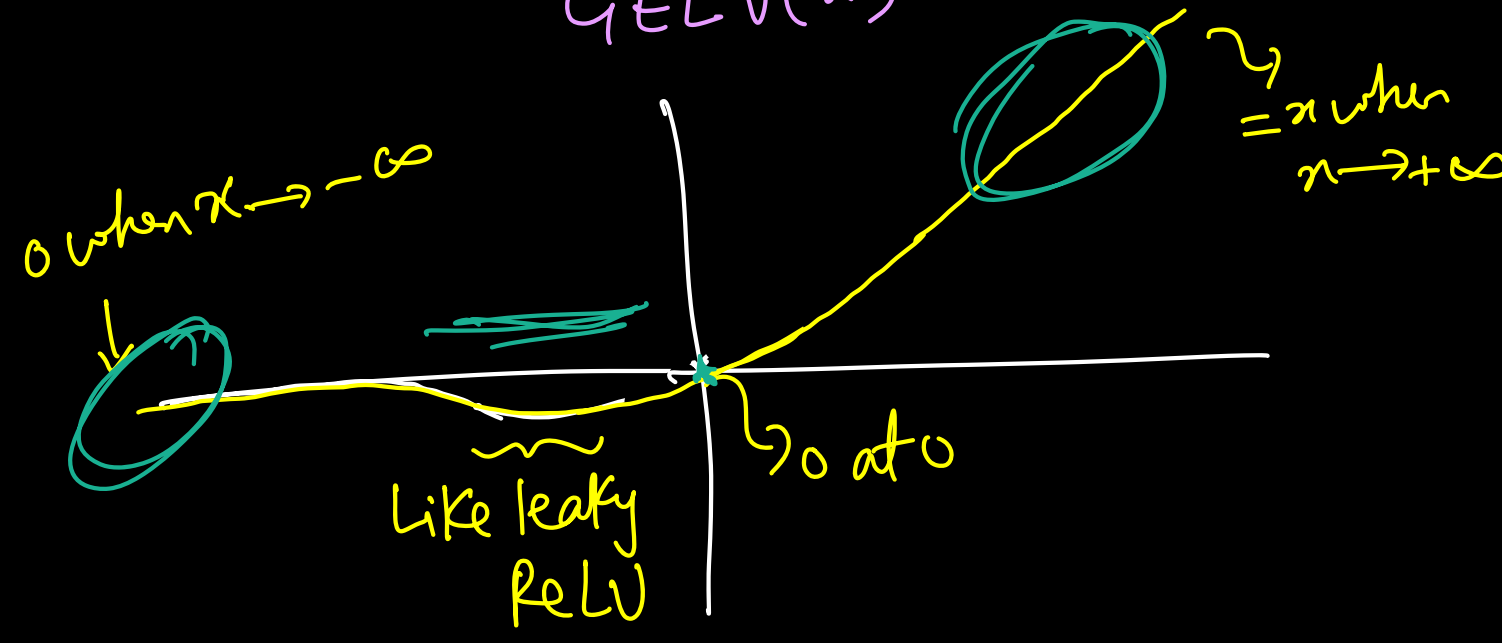
Gaussian distribution



$\Phi(x)$
CDF



GELU(x)



2. GELU Activation

$$F = \underbrace{\text{GELU}(H)}_{6 \times 3072} W_2 + \underbrace{b_2}_{6 \times 768}$$

6×768 (underlined)

3072×768 (circled)

Compression back (underlined)

2. Residual + Layer Norm

$$X^{(1)} = \text{Layer Norm} (\hat{X} + F)$$

6×768 (underlined)

Some dimension as input to the layer!

prev. Info. (underlined)

New Info.